

# Updating formulas and algorithms for computing entropy and Gini index on time-changing data streams

Blaž Sovdat\*

March 26, 2014

## Abstract

Despite growing interest in data stream mining the most successful incremental learners, such as VFDT, still use periodic recomputation to update attribute information gains and Gini indices. This note provides simple incremental formulas and algorithms for computing entropy and Gini index on time-changing data streams.

## 1 Introduction

Information-theoretic entropy was introduced by Shannon in the celebrated 1948 paper [Sha48] and has since found a vast number of applications [CT06]. In machine learning, information gain [Mit97], defined as the expected entropy reduction after splitting a leaf on a given attribute, is one of the most popular impurity measures for decision tree learning.

However, within the data stream mining world, there is a need for computationally cheap update formulas — an alternative being the complete and expensive recomputation — to compute the entropy as new examples come in, changing the sample distribution. Concrete example of such scenario is found in the incremental decision tree learner CVFDT [HSD01]. The same holds for Gini index, another popular impurity measure used for decision tree and IF-THEN rule learning.

Our main contribution are theorems 4 and 5 that give simple update formulas for entropy as new examples come in and as some of the sample counts change, respectively, and algorithms 3 and 4 that estimate “current” entropy of a data stream using our formulas with sliding windows and fading factors. We do the same for Gini index. We are not aware of existence of such formulas.

This note is organized as follows. In section 2 we derive incremental formulas for Gini index and then use these formulas (subsection 2.1) with sliding windows and fading factors to get algorithms for computing Gini index of time-changing data streams. In section 3 we give analogous formulas for entropy and use them (subsection 3.1) with sliding windows and fading factors to get algorithms for estimating entropy of a time-changing data streams. We conclude the note in section 4.

## 2 Incremental formulas for Gini index

Let  $\{x_i\}_{i=1}^n$  be a sample of positive real numbers and let  $S_n := x_1 + x_2 + \dots + x_n$  be sum of sample elements. Furthermore, let  $p_i := x_i/S_n$ .

---

\*Jožef Stefan Institute, Jamova 39, 1000 Ljubljana, Slovenia. Email: [blaz.sovdat@ijs.si](mailto:blaz.sovdat@ijs.si).

**Definition 1.** Gini index of  $\{x_i\}_{i=1}^n$  is defined as

$$G_n := 1 - \sum_{i=1}^n p_i^2 = 1 - \sum_{i=1}^n \left( \frac{x_i}{S_n} \right)^2 = 1 - \frac{1}{S_n^2} \sum_{i=1}^n x_i^2.$$

The next claim gives update formula when one of the sample elements increases by one.

**Claim 1.** Let  $\{x_i\}_{i=1}^n$  be sample of positive real numbers and let  $G_n$  be sample Gini index. Suppose  $x_i$  changes to  $x_i + 1$  and let  $G'_n$  denote Gini index of the new sample. We then have

$$G'_n = 1 - \frac{1}{(S_n + 1)^2} (S_n^2 (G_n - 1) + 2x_i + 1). \quad (1)$$

*Proof.* Plug the new value for  $x_i$  and do the algebra:

$$\begin{aligned} G'_n &= 1 - \frac{1}{S_n'^2} \sum_{i=1}^n x_i'^2 \\ &= 1 - \frac{1}{(S_n + 1)^2} (S_n^2 G_n + 2x_i + 1 - S_n^2) \\ &= 1 - \frac{1}{(S_n + 1)^2} (S_n^2 (G_n - 1) + 2x_i + 1). \end{aligned}$$

□

The next theorem generalizes the previous claim 1.

**Theorem 1.** Let  $\{x_i\}_{i=1}^n$  be sample of positive real numbers and let  $G_n$  be sample Gini index. Suppose that  $i$ th sample element  $x_i$  increases by  $r_i > 0$  for  $i \in I$ , where  $I$  is index set (indices of elements that change). Define  $r := r_1 + r_2 + \dots + r_n$  with  $r_i := 0$  for  $i \notin I$ . We then have

$$G'_n = 1 - \frac{1}{(S_n + r)^2} \left( S_n^2 (G_n - 1) + \sum_{i \in I} (2x_i r_i + r_i^2) \right). \quad (2)$$

*Proof.* Similarly as before, we do the algebra to get the result:

$$\begin{aligned} G'_n &= 1 - \frac{1}{(S_n + r)^2} \sum_{i=1}^n (x_i + r_i)^2 \\ &= 1 - \frac{1}{(S_n + r)^2} \sum_{i=1}^n (x_i^2 + 2x_i r_i + r_i^2) \\ &= 1 - \frac{1}{(S_n + r)^2} \left( \sum_{i=1}^n x_i^2 + \sum_{i \in I} (2x_i r_i + r_i^2) \right) \\ &= 1 - \frac{1}{(S_n + r)^2} \left( S_n^2 - S_n^2 \sum_{i=1}^n \frac{x_i^2}{S_n^2} - S_n^2 + \sum_{i \in I} (2x_i r_i + r_i^2) \right) \\ &= 1 - \frac{1}{(S_n + r)^2} \left( S_n^2 (G_n - 1) + \sum_{i \in I} (2x_i r_i + r_i^2) \right). \end{aligned}$$

□

Note that the number of required operations in (2) grows linearly with the number of elements that changed.

The next claim gives update formula for Gini index when new element enters the sample.

**Claim 2.** *Let  $\{x_i\}_{i=1}^n$  be sample of positive real numbers and let  $S_n$  and  $G_n$  be the sum of sample elements and sample Gini index, respectively. Suppose new element,  $x_{n+1}$ , enters the sample. Gini index then becomes*

$$G_{n+1} = 1 - \frac{1}{(S_n + x_{n+1})^2} (S_n^2(G_n - 1) + x_{n+1}^2). \quad (3)$$

*Proof.* We clearly have

$$\begin{aligned} G_{n+1} &= 1 - \frac{1}{(S_n + x_{n+1})^2} \left( \sum_{i=1}^n x_i^2 + x_{n+1}^2 \right) \\ &= 1 - \frac{1}{(S_n + x_{n+1})^2} \left( S_n^2 - S_n^2 \sum_{i=1}^n \frac{x_i^2}{S_n^2} - S_n^2 + x_{n+1}^2 \right) \\ &= 1 - \frac{1}{(S_n + x_{n+1})^2} (S_n^2(G_n - 1) + x_{n+1}^2). \end{aligned}$$

□

The next theorem generalizes claim 2, giving update formula for Gini index when we “concatenate” two samples.

**Theorem 2.** *Let  $\{x_i\}_{i=1}^n$  and  $\{y_i\}_{i=1}^m$  be samples of positive real numbers and let  $S_n$  and  $S_m$  be sums of sample elements. Furthermore, let  $G_n$  and  $G_m$  be sample Gini indices and define*

$$z_i := \begin{cases} x_i, & 1 \leq i \leq n, \\ y_{i-n}, & n+1 \leq i \leq n+m. \end{cases}$$

*Then the Gini index becomes*

$$G_{n+m} = 1 - \frac{1}{(S_n + S_m)^2} (S_n^2(G_n - 1) + S_m^2(G_m - 1)).$$

*Proof.* By definition we have

$$\begin{aligned} G_{n+m} &= 1 - \frac{1}{(S_n + S_m)^2} \sum_{i=1}^{n+m} z_i^2 \\ &= 1 - \frac{1}{(S_n + S_m)^2} \left( \sum_{i=1}^n x_i^2 + \sum_{i=1}^m y_i^2 \right) \\ &= 1 - \frac{1}{(S_n + S_m)^2} \left( S_n^2 \sum_{i=1}^n \frac{x_i^2}{S_n^2} + S_m^2 \sum_{i=1}^m \frac{y_i^2}{S_m^2} \right) \\ &= 1 - \frac{1}{(S_n + S_m)^2} (S_n^2(G_n - 1) + S_m^2(G_m - 1)). \end{aligned}$$

□

**Theorem 3.** Let  $\{x_i\}_{i=1}^n$  and  $\{y_i\}_{i=1}^n$  be sample of positive real numbers, let  $R_n$  in  $S_n$  be sums of sample elements, and let  $F_n$  and  $G_n$  be sample Gini indices. Furthermore, let  $z_i := x_i + y_i$ . Then Gini index of the “union” becomes

$$G = 1 - \frac{1}{(R_n + S_n)^2} \left( R_n^2 (F_n - 1) + S_n^2 (G_n - 1) + 2 \sum_{i=1}^n x_i y_i \right)$$

*Proof.* By definition we have

$$F_n = 1 - \frac{1}{R_n^2} \sum_{i=1}^n x_i^2$$

and

$$G_n = 1 - \frac{1}{S_n^2} \sum_{i=1}^n y_i^2.$$

First part of the formula follows by theorem 1, while the last part,  $2(x_1 y_1 + x_2 y_2 + \dots + x_n y_n)$ , is not stored anywhere and has to be recomputed.  $\square$

## 2.1 Algorithms for computing Gini index on time-changing data streams

Data streams are inherently changing and we are usually interested in “recent” Gini index. In this section we propose two algorithms for this problem — algorithm 1 uses sliding windows, while algorithm 2 uses fading factors to capture “recent”.

Algorithm 1 computes Gini index of the last  $w \in \mathbb{N}$  stream elements. It achieves this using sliding window of size  $w$ , meaning its space complexity is  $O(w)$ . Note that  $w$  is user-defined parameter, which indicates what subset of stream elements is “recent”.

Algorithm 2 computes “recent” Gini index using fading factors — element contributions are weighted with  $\{\alpha^k : k \in \mathbb{N}\}$  for some fixed  $\alpha \in (0, 1]$  according to element’s “age”. Note that the fading factor  $\alpha$  defines what recent means and that this algorithm has small constant space complexity.

## 3 Incremental formulas for entropy

In this section we give simple incremental formulas for entropy.

Let  $\{x_i\}_{i=1}^n$  be sample of positive real numbers and let  $S_n := x_1 + x_2 + \dots + x_n$  be sum of sample elements. Define entropy as

$$H_n := - \sum_{i=1}^n \frac{x_i}{S_n} \log_2 \frac{x_i}{S_n}.$$

We first prove the following technical lemma.

**Lemma 1.** Let  $\{x_i\}_{i=1}^n$  be sample of positive real numbers and let  $H_n$  and  $S_n$  be sample entropy and sum of sample elements, respectively. For any positive real number  $R > 0$  we have

$$- \sum_{i=1}^n \frac{x_i}{S_n + R} \log_2 \frac{x_i}{S_n + R} = \frac{S_n}{S_n + R} H_n - \frac{S_n}{S_n + R} \log_2 \frac{S_n}{S_n + R}. \quad (4)$$

---

**Algorithm 1** Computing Gini index using sliding windows.

---

```
1: Let  $W := \{\}$  be sliding window
2: Let  $n := 0$  be the number of all examples and let in  $n_i := 0$  be the number of examples from  $i$ th class
3: Let  $g := 0$  be the current Gini index
4: for  $x \in S$  do ▷ TODO: Add usage example for ADD and DEL
5:   if  $|W| > w$  then
6:     Remove the oldest element  $x'$ , labeled with  $i$ th class, from sliding window  $W$ 
7:     Update  $g := \text{DEC}(g, n, n_i)$ 
8:     Add  $W := W \cup \{x\}$  element labeled with  $i$ th class
9:     Update  $g := \text{INC}(g, n, n_i)$ 
10: function  $\text{ADD}(g, n, n_i)$  ▷ Append
11:   Update  $n := n + n_i$ 
12:   return  $1 - \frac{1}{n^2} ((n - n_i)^2(g - 1) + n_i^2)$ 
13: function  $\text{DEL}(g, n, n_i)$  ▷ Delete
14:   Update  $n := n - n_i$ 
15:   return  $1 - \frac{1}{(n + n_i)^2} ((g - 1)n^2 + n_i^2)$ 
16: function  $\text{INC}(g, n, n_i)$  ▷ Increment
17:   Update  $n := n + 1$  in  $n_i := n_i + 1$ 
18:   return  $1 - \frac{1}{n^2} ((n - 1)^2(g - 1) + 2n_i - 1)$ 
19: function  $\text{DEC}(g, n, n_i)$  ▷ Decrement
20:   Update  $n := n - 1$  in  $n_i := n_i - 1$ 
21:   return  $\frac{1}{n^2} ((n + 1)^2(1 - g) - 2n_i - 1)$ 
```

---

---

**Algorithm 2** Computing Gini index using fading factors.

---

**Require:** Fading factor  $\alpha \in (0, 1]$  and data stream  $S$ .

**Ensure:** Recent Gini index anytime.

```
1: Let  $n := 0$  be the number of all examples and let  $n_i := 0$  be the number of examples from  $i$ th class
2: Let  $g := 0$  be the current Gini index
3: for  $x \in S$  do
4:   Update Gini index  $g := 1 - \frac{1}{(n + 1)^2} (n^2(\alpha g - 1) + 2n_i + 1)$ 
5:   Update counts  $n := n + 1$  and  $n_i := n_i + 1$ 
```

---

*Proof.* Write  $\frac{x_i}{R+S_n} = 1 \cdot \frac{x_i}{R+S_n} = \frac{S_n}{S_n} \frac{x_i}{R+S_n} = \frac{x_i}{S_n} \frac{S_n}{R+S_n}$ . We then clearly have

$$\begin{aligned}
-\sum_{i=1}^n \frac{x_i}{S_n+R} \log_2 \frac{x_i}{S_n+R} &= -\sum_{i=1}^n \frac{x_i}{S_n} \frac{S_n}{S_n+R} \log_2 \left( \frac{x_i}{S_n} \frac{S_n}{S_n+R} \right) \\
&= -\sum_{i=1}^n \frac{x_i}{S_n} \frac{S_n}{S_n+R} \left( \log_2 \frac{x_i}{S_n} + \log_2 \frac{S_n}{S_n+R} \right) \\
&= -\frac{S_n}{S_n+R} \sum_{i=1}^n \frac{x_i}{S_n} \log_2 \frac{x_i}{S_n} - \frac{S_n}{S_n+R} \sum_{i=1}^n \frac{x_i}{S_n} \log_2 \frac{S_n}{S_n+R} \\
&= \frac{S_n}{S_n+R} H_n - \frac{S_n}{S_n+R} \log_2 \frac{S_n}{S_n+R} \sum_{i=1}^n \frac{x_i}{S_n} \\
&= \frac{S_n}{S_n+R} H_n - \frac{S_n}{S_n+R} \log_2 \frac{S_n}{S_n+R}.
\end{aligned}$$

□

The next claim gives simple update formula when a new positive real number  $x_i > 0$  enters the sample.

**Claim 3** ([Sch13]). *Let  $H_n$  and  $S_n$  be sample entropy and sum of sample elements and suppose that a new positive real number  $x_{n+1} > 0$  enters the sample. We then have*

$$H_{n+1} = \frac{S_n}{S_{n+1}} H_n - \frac{S_n}{S_{n+1}} \log_2 \frac{S_n}{S_{n+1}} - \frac{x_{n+1}}{S_{n+1}} \log_2 \frac{x_{n+1}}{S_{n+1}}. \quad (5)$$

*Proof.* By definition we have

$$\begin{aligned}
H_{n+1} &= -\sum_{i=1}^{n+1} \frac{x_i}{S_{n+1}} \log_2 \frac{x_i}{S_{n+1}} \\
&= -\frac{x_{n+1}}{S_{n+1}} \log_2 \frac{x_{n+1}}{S_{n+1}} - \sum_{i=1}^n \frac{x_i}{S_{n+1}} \log_2 \frac{x_i}{S_{n+1}} \\
&= \frac{S_n}{S_{n+1}} H_n - \frac{S_n}{S_{n+1}} \log_2 \frac{S_n}{S_{n+1}} - \frac{x_{n+1}}{S_{n+1}} \log_2 \frac{x_{n+1}}{S_{n+1}},
\end{aligned}$$

with the last equality following from lemma 1. □

The next theorem generalizes the claim and gives formula for entropy of “concatenation” of samples, given sample entropies  $H_n$  and  $H_m$  and sums of sample elements  $S_n$  and  $S_m$ .

**Theorem 4.** *Let  $\{x_i\}_{i=1}^n$  in  $\{y_i\}_{i=1}^m$  be samples of positive real numbers and let  $S_n := x_1 + x_2 + \dots + x_n$  in  $R_m := y_1 + y_2 + \dots + y_m$  be sums of sample elements. Furthermore let  $H_n$  and  $H_m$  be sample entropies. Define*

$$z_i := \begin{cases} x_i, & 1 \leq i \leq n, \\ y_{i-n}, & n+1 \leq i \leq n+m, \end{cases}$$

and let  $Z_{n+m} := z_1 + z_2 + \dots + z_{n+m} = S_n + R_m$ . We then have

$$H_{n+m} = \frac{S_n}{Z_{n+m}} H_n - \frac{S_n}{Z_{n+m}} \log_2 \frac{S_n}{Z_{n+m}} + \frac{R_m}{Z_{n+m}} H_m - \frac{R_m}{Z_{n+m}} \log_2 \frac{R_m}{Z_{n+m}}. \quad (6)$$

*Proof.* Similarly as before, we have

$$\begin{aligned} -\sum_{i=1}^{n+m} \frac{z_i}{Z_{n+m}} \log_2 \frac{z_i}{Z_{n+m}} &= -\sum_{i=1}^n \frac{x_i}{Z_{n+m}} \log_2 \frac{x_i}{Z_{n+m}} - \sum_{i=1}^m \frac{y_i}{Z_{n+m}} \log_2 \frac{y_i}{Z_{n+m}} \\ &= \frac{S_n}{Z_{n+m}} H_n - \frac{S_n}{Z_{n+m}} \log_2 \frac{S_n}{Z_{n+m}} - \frac{R_m}{Z_{n+m}} H_m - \frac{R_m}{Z_{n+m}} \log_2 \frac{R_m}{Z_{n+m}}, \end{aligned}$$

where the last equality follows by applying lemma 1 twice.  $\square$

Note that claim 3 is a corollary of theorem 4, if apply equation (6) to  $H_n$  and  $x_{n+1}$  and think of  $x_{n+1}$  as a sample with a single element.

Theorem 5 gives formula for entropy when some of the elements  $x_{i_j}$  for  $1 \leq j \leq k < n$  increase.

**Theorem 5.** *Let  $\{x_i\}_{i=1}^n$  be a sample of positive real numbers and let  $S_n$  be sum of sample elements. Let  $H_n$  be sample entropy. Suppose  $x_{i_j}$  increases by  $r_j > 0$  for  $1 \leq j \leq k < n$  and let  $r := r_1 + r_2 + \dots + r_k$ . Then the entropy  $H_n$  becomes*

$$\frac{S_n}{S_n + r} H_n - \frac{S_n}{S_n + r} \log_2 \frac{S_n}{S_n + r} - \sum_{j=1}^k \left( \frac{x_{i_j} + r_j}{S_n + r} \log_2 \frac{x_{i_j} + r_j}{S_n + r} - \frac{x_{i_j}}{S_n} \log_2 \frac{x_{i_j}}{S_n} \right). \quad (7)$$

*Proof.* The idea is to think of  $\frac{x_{i_j} + r_j}{S_n + r}$  as new elements, apply theorem 4, and subtract “old” elements  $\frac{x_{i_j}}{S_n}$ .  $\square$

Note that the number of required operations grow linearly with the number of changed elements — if  $k$  elements change, we only need  $O(k)$  operations. (Also note that formulas become (numerically) problematic when  $x_n$  is small compared to  $S_n$ .<sup>1</sup>)

### 3.1 Algorithms for computing entropy on time-changing data streams

We now give algorithms for computing “recent” entropy — algorithm 3 uses sliding windows, while algorithm 4 uses fading factors.

Algorithm 3, similarly as its Gini-index-analog, algorithm 1, sliding window size  $w \in \mathbb{N}$ , which defines what subset of stream elements is recent. The space complexity is clearly  $O(w)$ .

Ideally we would want an algorithm that adapts sliding window size — similarly as ADWIN [BG07] does — because  $w$  changes with time due to time-changing nature of data streams.

Algorithm 4, similarly as algorithm 2, defines “recent” using fading factors  $\alpha \in (0, 1]$ . Element contributions are weighted with  $\{\alpha^k : k \in \mathbb{N}\}$  according to element “age”. The algorithm has small constant space complexity.

## 4 Conclusion

We derived simple incremental formulas and algorithms for computing entropy and Gini index from time-changing data streams. In future work we plan to address numerical stability issues of the derived formulas.

---

<sup>1</sup>Thanks to Andrej Muhič for providing pathological scenario.

---

**Algorithm 3** Computing entropy using sliding windows.

---

**Require:** Sliding window size  $w \in \mathbb{N}$ .

**Ensure:** Current entropy.

```

1: Let  $W := \{\}$  be sliding window and let  $h := 0$  be the current entropy.
2: Let  $n := 0$  be the number of all examples and let  $n_i := 0$  be the number of examples in  $i$ th class.
3: for  $x \in S$  do ▷ Missing usage example for ADD and DEC.
4:   if  $|W| > w$  then
5:     Remove the oldest element  $x'$  from sliding window  $W$ 
6:     Update  $h := \text{DEC}(h, n, x')$ 
7:   Add  $W := W \cup \{x\}$ 
8:   Update  $h := \text{INC}(h, n, x)$ 
9: function  $\text{ADD}(h, n, n_i)$ 
10:   Update  $n := n + n_i$ 
11:   return  $\frac{n}{n - n_i} h - \frac{n}{n - n_i} \log_2 \frac{n}{n - n_i} - \frac{n_i}{n} \log_2 \frac{n_i}{n}$ 
12: function  $\text{DEL}(h, n, n_i)$ 
13:   Update  $n := n - n_i$ 
14:   return  $\frac{n}{n + n_i} \left( h + \frac{n + n_i}{n} \log_2 \frac{n + n_i}{n} + \frac{n_i}{n + n_i} \log_2 \frac{n_i}{n + n_i} \right)$ 
15: function  $\text{INC}(h, n, n_i)$ 
16:   Update  $n := n + 1$  in  $n_i := n_i + 1$ 
17:   return  $\frac{n}{n - 1} h - \frac{n}{n - 1} \log_2 \frac{n}{n - 1} - \frac{n_i}{n - 1} \log_2 \frac{n_i}{n - 1}$ 
18: function  $\text{DEC}(h, n, n_i)$ 
19:   Update  $n := n - 1$  in  $n_i := n_i - 1$ 
20:   return  $\frac{n + 1}{n} \left( h + \frac{n}{n + 1} \log_2 \frac{n}{n + 1} + \frac{n_i + 1}{n + 1} \log_2 \frac{n_i + 1}{n + 1} \right)$ 

```

---

## Acknowledgments

The author is grateful to Marko Robnik-Šikonja, Martin Vuk, and Zoran Bosnić from Faculty of Computer and Information Science at University of Ljubljana for useful comments and suggestions and to Blaž Fortuna, Andrej Muhič, Jan Rupnik, Marko Grobelnik and others from Artificial Intelligence Laboratory at Jožef Stefan Institute.

## References

- [BG07] Albert Bifet and Ricard Gavalda. Learning from time-changing data with adaptive windowing. In *SDM*, volume 7, page 2007. SIAM, 2007.
- [CT06] Thomas M. Cover and Joy A. Thomas. *Elements of Information Theory*. Wiley-Interscience, 2nd edition, 2006.
- [HSD01] Geoff Hulten, Laurie Spencer, and Pedro Domingos. Mining time-changing data streams. In *Proceedings of the seventh ACM SIGKDD international conference on Knowledge discovery and data mining*, KDD '01, pages 97–106, New York, NY, USA, 2001. ACM.
- [Mit97] Thomas M. Mitchell. *Machine Learning*. McGraw-Hill, Inc., New York, NY, USA, 1st edition, 1997.



---

**Algorithm 4** Computing entropy from time-changing data streams.

---

**Require:** Fading factor  $\alpha \in (0, 1]$ .**Ensure:** Current Gini index anytime.

- 1: Let  $n := 0$  be the number of all examples and let  $n_i := 0$  be the number of example in  $i$ th class
  - 2: Let  $h := 0$  be the current entropy
  - 3: **for**  $x \in S$  **do**
  - 4:   Let  $i$  be class of element  $x$
  - 5:   Update entropy  $h := \alpha \frac{n}{n+1} h - \frac{n+1}{n} \log_2 \frac{n+1}{n} - \frac{n_i+1}{n+1} \log_2 \frac{n_i+1}{n+1}$
  - 6:   Update counts  $n := n+1$  and  $n_i := n_i+1$
- 

[Sch13] André Schlichting. Incremental entropy computation. MathOverflow, 2013. [Online; accessed 2013-06-17].

[Sha48] Claude E. Shannon. A mathematical theory of communication. *The Bell System Technical Journal*, 27:379–423, 623–656, July, October 1948.